

RESEARCH INITIATIVE FOR SECURE AI-POWERED DECISION-MAKING

PROBLEM

Competitive pressures force more and more outsourcing of decision-making to AI. However, it is unclear whether the models inherit or resist human reactions to “**judgment-impairing conditions**” such as time-pressure, information ambiguity, high-stakes situations, and more.

We need a better understanding of how models maintain a standard for reasoning and how humans can engage with their output under such conditions. Otherwise, humanity risks gradually losing control over increasingly AI-powered decision-making.

RESEARCH AIM

Uncover harmful weaknesses in decision-making between people and AI in situations where judgement is impaired.

RESEARCH QUESTIONS

Technical: How can we measure the robustness of model performance under judgment-impairing conditions?

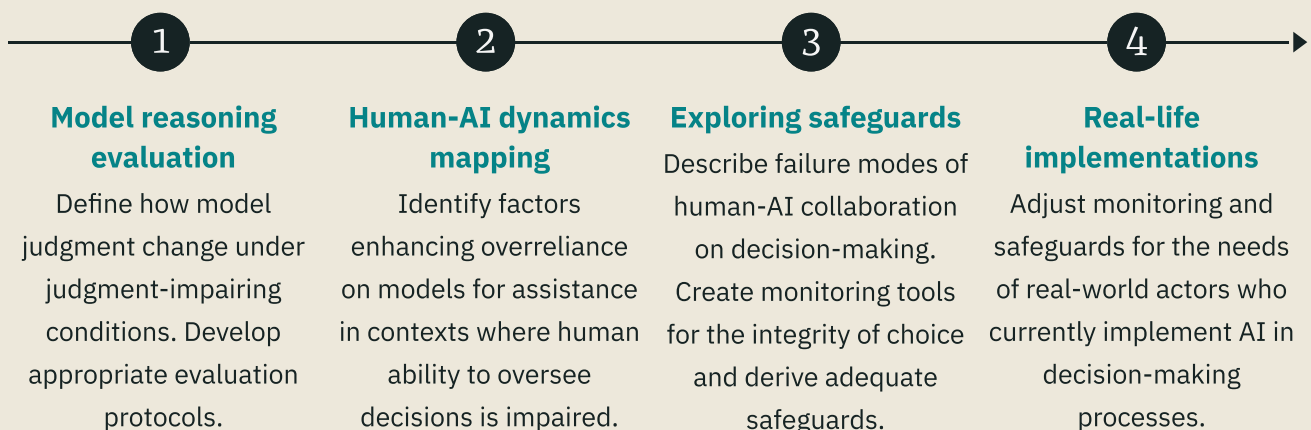
Cognitive: How do such conditions impact human decision-makers’ overreliance and how transparent is this risk to them?

Remedial: How can the identified vulnerabilities be used to shape model evaluations, early-warning alerts, and safeguards?

IMPACT

A foundation for monitoring the integrity of human choice and enable the development of safeguards against disempowerment in areas facing the most immediate risk.

STAGES ///



Does this sound interesting? Reach out. Seeking both researchers and operations talent.